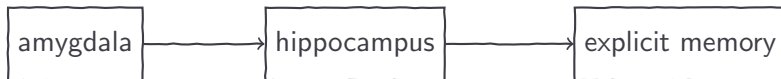MAX-PLANCK-GESELLSCHAFT

# Causal interpretation rules for encoding and decoding models in neuroimaging

Sebastian Weichwald, Timm Meyer, Ozan Özdenizci[§],
Bernhard Schölkopf, Tonio Ball[‡], Moritz Grosse-Wentrup

MPI for Intelligent Systems, [§]Sabanci University, [‡]University of Freiburg

# Motivation

*Hippocampal activity in this study was correlated with amygdala activity, supporting the view that the amygdala* **enhances** *explicit memory by* **modulating** *activity in the hippocampus.*

(S. Hamann, *Trends in Cognitive Sciences*, 2001)

*We tested [...] whether pre-stimulus alpha oscillations measured with electroencephalography (EEG) **influence** the encoding of items into working memory.*

(Myers et al., *Journal of Neuroscience*, 2014)

# Approach

5

# Encoding and decoding models in neuroimaging

Trial 3    Trial 4    Trial 5    Trial 6    Trial 7

L          L          R          R          R

encoding

decoding

e.g. mean difference
between conditions

e.g. classifier for
experimental conditions

"Feature shows significant variation across experimental conditions?"

"Feature helpful for predicting the experimental condition?"

"Feature shows significant variation across experimental conditions?"

"Feature helpful for predicting the experimental condition?"

?
relevant feature ↭ cognitive process

# Causal interpretation of encoding and decoding models

Let's set out the causal component of already performed analyses..

Let's set out the causal component of already performed analyses..

stimulus- vs response-based

feature relevance ↭ marginal/conditional dependence

⤳ 16 causal interpretation rules



9

Let's set out the causal component of already performed analyses..

stimulus- vs response-based

feature relevance $\leftrightarrow$ marginal/conditional dependence

$\rightsquigarrow$ 16 causal interpretation rules
simple

$S$ — $\vec{X} = \{X_1, ..., X_d\}$ → $R$

stimulus      brain state features      response

10

| | stimulus-based | | response-based |
|---|---|---|---|
| $p(\vec{X}|S)$ | encoding | | $p(\vec{X}|R)$ |
| $p(S|\vec{X})$ | decoding | | $p(R|\vec{X})$ |

|  | stimulus-based |  | response-based |  |
|---|---|---|---|---|
| $p(\vec{X}|S)$ | causal | encoding |  | $p(\vec{X}|R)$ |
| $p(S|\vec{X})$ |  | decoding | causal | $p(R|\vec{X})$ |

|  | stimulus-based |  | response-based |  |
|---|---|---|---|---|
| $p(\vec{X}|S)$ | causal | encoding | *anti*-causal | $p(\vec{X}|R)$ |
| $p(S|\vec{X})$ | *anti*-causal | decoding | causal | $p(R|\vec{X})$ |

$$p(X_i | C = c_1) \overset{?}{\neq} p(X_i | C = c_2)$$

$$p(X_i | C = c_1) \overset{?}{\neq} p(X_i | C = c_2) \qquad X_i \not\perp C$$

$$p(X_i | C = c_1) \overset{?}{\neq} p(X_i | C = c_2) \qquad\qquad X_i \not\perp C$$

$$p(C | \vec{X}) \overset{?}{\neq} p(C | \vec{X} \smallsetminus X_i)$$

$$p(X_i | C = c_1) \overset{?}{\neq} p(X_i | C = c_2) \qquad\qquad X_i \not\perp C$$

$$p(C | \vec{X}) \overset{?}{\neq} p(C | \vec{X} \smallsetminus X_i) \qquad\qquad X_i \not\perp C | \vec{X} \smallsetminus X_i$$

12

| | Feature $X_i$ relevant? | | |
| | Encoding | Decoding | Causal interpretation |
| --- | --- | --- | --- |
| Stimulus-based | $\times$ $\checkmark$ | $\times$ $\checkmark$ | |
| Response-based | $\times$ $\checkmark$ | $\times$ $\checkmark$ | |

12

| | Feature $X_i$ relevant? | | |
|---|---|---|---|
| | Encoding | Decoding | Causal interpretation |
| Stimulus-based | $\times$ | | |
| | $\checkmark$ | $X_i \not\perp R$ | |
| | | $\times$ | |
| | | $X_i \leftarrow h \rightarrow R$ | |
| | | $\checkmark$ | |
| Response-based | $\times$ | $X_i \rightarrow R$ | |
| | $\checkmark$ | | |
| | | $\times$ | |
| | | $\checkmark$ | |

|  | Feature $X_i$ relevant? | | Causal interpretation |
| --- | --- | --- | --- |
|  | Encoding | Decoding |  |
| Stimulus-based | × | | |
| | √ | | |
| | | × | |
| | | √ | |
| Response-based | × | | |
| | √ | | |
| | | × | |
| | | √ | |

12

| | Feature $X_i$ relevant? | | |
| | Encoding | Decoding | Causal interpretation |
| --- | --- | --- | --- |
| Stimulus-based | $\times$ | | |
| | $\sqrt{}$ | | |
| | | $\times$ | |
| | | $\sqrt{}$ | |
| Response-based | $\times$ | | |
| | $\sqrt{}$ | | inconclusive |
| | | $\times$ | |
| | | $\sqrt{}$ | |

12

| | Feature $X_i$ relevant? | | |
| | Encoding | Decoding | Causal interpretation |
|---|---|---|---|
| **Stimulus-based** | $\times$ | | |
| | $\checkmark$ | | |
| | | $\times$ | inconclusive |
| | | $\checkmark$ | inconclusive |
| **Response-based** | $\times$ | | |
| | $\checkmark$ | | inconclusive |
| | | $\times$ | inconclusive |
| | | $\checkmark$ | inconclusive |

12

| | Feature $X_i$ relevant? | | |
| | Encoding | Decoding | Causal interpretation |
| --- | --- | --- | --- |
| Stimulus-based | × | | no effect of $S$ |
| | √ | | effect of $S$ |
| | | × | inconclusive |
| | | √ | inconclusive |
| Response-based | × | | no cause of $R$ |
| | √ | | inconclusive |
| | | × | inconclusive |
| | | √ | inconclusive |

12

| | Feature $X_i$ relevant? | | Causal interpretation |
| | Encoding | Decoding | |
|---|---|---|---|
| **Stimulus-based** | $\checkmark$ | $\checkmark$ | |
| | $\checkmark$ | $\times$ | |
| | $\times$ | $\checkmark$ | |
| | $\times$ | $\times$ | |
| **Response-based** | $\checkmark$ | $\checkmark$ | |
| | $\checkmark$ | $\times$ | |
| | $\times$ | $\checkmark$ | |
| | $\times$ | $\times$ | |

13

|  | Feature $X_i$ relevant? | | Causal interpretation |
|  | Encoding | Decoding |  |
| --- | --- | --- | --- |
| **Stimulus-based** | $\checkmark$ | $\checkmark$ | |
|  | $\checkmark$ | $\times$ | |
|  | $\times$ | $\checkmark$ | |
|  | $\times$ | $\times$ | |
| **Response-based** | $\checkmark$ | $\checkmark$ | inconclusive |
|  | $\checkmark$ | $\times$ | |
|  | $\times$ | $\checkmark$ | |
|  | $\times$ | $\times$ | |

| | Feature $X_i$ relevant? | | |
| --- | --- | --- | --- |
| | Encoding | Decoding | Causal interpretation |
| Stimulus-based | $\sqrt{}$ | $\sqrt{}$ | |
| | $\sqrt{}$ | $\times$ | |
| | $\times$ | $\sqrt{}$ | |
| | $\times$ | $\times$ | |
| Response-based | $\sqrt{}$ | $\sqrt{}$ | inconclusive |
| | $\sqrt{}$ | $\times$ | |
| | $\times$ | $\sqrt{}$ | |
| | $\times$ | $\times$ | |

$X_i \not\perp\!\!\!\perp S$  and  $X_i \perp\!\!\!\perp S | \vec{X} \smallsetminus X_i$

$S \dashrightarrow X_i$  indirectly

| | Feature $X_i$ relevant? | | |
| | Encoding | Decoding | Causal interpretation |
|---|---|---|---|
| Stimulus-based | $\checkmark$ | $\checkmark$ | |
| | $\checkmark$ | $\times$ | |
| | $\times$ | $\checkmark$ | |
| | $\times$ | $\times$ | |
| Response-based | $\checkmark$ | $\checkmark$ | inconclusive |
| | $\checkmark$ | $\times$ | |
| | $\times$ | $\checkmark$ | |
| | $\times$ | $\times$ | |

13

| | Feature $X_i$ relevant? | | |
| --- | --- | --- | --- |
| | Encoding | Decoding | Causal interpretation |
| **Stimulus-based** | $\checkmark$ | $\checkmark$ | |
| | $\checkmark$ | $\times$ | indirect effect of $S$ |
| | $\times$ | $\checkmark$ | |
| | $\times$ | $\times$ | |
| **Response-based** | $\checkmark$ | $\checkmark$ | inconclusive |
| | $\checkmark$ | $\times$ | |
| | $\times$ | $\checkmark$ | |
| | $\times$ | $\times$ | |

13

| | Feature $X_i$ relevant? | | |
| | Encoding | Decoding | Causal interpretation |
|---|---|---|---|
| Stimulus-based | $\checkmark$ | $\checkmark$ | effect of $S$ |
| | $\checkmark$ | $\times$ | indirect effect of $S$ |
| | $\times$ | $\checkmark$ | provides context |
| | $\times$ | $\times$ | no effect of $S$ |
| Response-based | $\checkmark$ | $\checkmark$ | inconclusive |
| | $\checkmark$ | $\times$ | no direct cause of $R$ |
| | $\times$ | $\checkmark$ | provides context |
| | $\times$ | $\times$ | no cause of $R$ |

13

# Empirical example

$\alpha_{\mathsf{IC}_1}$  $\alpha_{\mathsf{IC}_2}$  $\alpha_{\mathsf{IC}_3}$  $\alpha_{\mathsf{IC}_4}$  $\alpha_{\mathsf{IC}_5}$  $\alpha_{\mathsf{IC}_6}$

$\alpha_{\mathsf{IC}_1}$  $\alpha_{\mathsf{IC}_2}$  $\alpha_{\mathsf{IC}_3}$  $\alpha_{\mathsf{IC}_4}$  $\alpha_{\mathsf{IC}_5}$  $\alpha_{\mathsf{IC}_6}$

*p*-values

Encoding

Decoding

| *p*-values | $\alpha_{IC_1}$ | $\alpha_{IC_2}$ | $\alpha_{IC_3}$ | $\alpha_{IC_4}$ | $\alpha_{IC_5}$ | $\alpha_{IC_6}$ |
|---|---|---|---|---|---|---|
| Encoding | 0 | 0 | 0 | 0 | 0 | 0 |
| Decoding | 0 | 0 | 0.50 | 0.34 | 0.79 | 0.13 |

| $p$-values | $\alpha_{\mathsf{IC}_1}$ | $\alpha_{\mathsf{IC}_2}$ | $\alpha_{\mathsf{IC}_3}$ | $\alpha_{\mathsf{IC}_4}$ | $\alpha_{\mathsf{IC}_5}$ | $\alpha_{\mathsf{IC}_6}$ |
|---|---|---|---|---|---|---|
| Encoding | 0 | 0 | 0 | 0 | 0 | 0 |
| Decoding | 0 | 0 | 0.50 | 0.34 | 0.79 | 0.13 |

| $p$-values | $\alpha_{IC_1}$ | $\alpha_{IC_2}$ | $\alpha_{IC_3}$ | $\alpha_{IC_4}$ | $\alpha_{IC_5}$ | $\alpha_{IC_6}$ |
|---|---|---|---|---|---|---|
| Encoding | 0 | 0 | 0 | 0 | 0 | 0 |
| Decoding | 0 | 0 | 0.50 | 0.34 | 0.79 | 0.13 |

| $p$-values | $\alpha_{\mathsf{IC}_1}$ | $\alpha_{\mathsf{IC}_2}$ | $\alpha_{\mathsf{IC}_3}$ | $\alpha_{\mathsf{IC}_4}$ | $\alpha_{\mathsf{IC}_5}$ | $\alpha_{\mathsf{IC}_6}$ |
|---|---|---|---|---|---|---|
| Encoding | 0 | 0 | 0 | 0 | 0 | 0 |
| Decoding | 0 | 0 | 0.50 | 0.34 | 0.79 | 0.13 |

- ‣ instruction to plan a reaching movement is causal for all $\alpha_{\mathsf{IC}_i}$
- ‣ $\alpha_{\mathsf{IC}_3}, ..., \alpha_{\mathsf{IC}_6}$ are only indirect effects

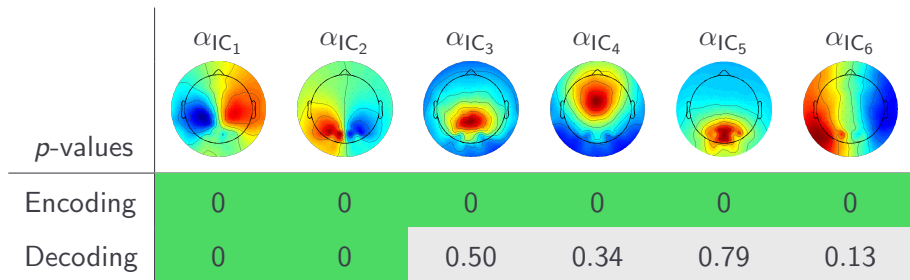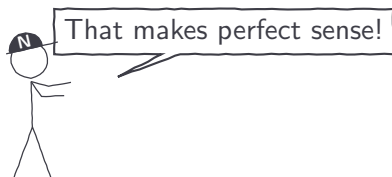| $p$-values | $\alpha_{IC_1}$ | $\alpha_{IC_2}$ | $\alpha_{IC_3}$ | $\alpha_{IC_4}$ | $\alpha_{IC_5}$ | $\alpha_{IC_6}$ |
|---|---|---|---|---|---|---|
| Encoding | 0 | 0 | 0 | 0 | 0 | 0 |
| Decoding | 0 | 0 | 0.50 | 0.34 | 0.79 | 0.13 |

‣ instruction to plan a reaching movement is causal for all $\alpha_{IC_i}$

‣ $\alpha_{IC_3}, ..., \alpha_{IC_6}$ are only indirect effects

That makes perfect sense!

# Wrap-up

feature relevance

feature relevance $\swarrow^\nearrow$ (conditional) (in)dependence

causal structure

feature relevance ⤴ (conditional) (in)dependence ⤴

feature relevance $\nearrow$ (conditional) (in)dependence $\nearrow$ causal structure

- ‣ simple interpretation rules
- ‣ reinterpretation of previous results?
- ‣ resolve recently discussed issues

Sebastian Weichwald, Timm Meyer, Ozan Özdenizci, Bernhard Schölkopf, Tonio Ball,
Moritz Grosse-Wentrup:

▸ Causal interpretation rules for encoding and decoding models in
  neuroimaging. *NeuroImage*, 2015.

▸ Causal and anti-causal learning in pattern recognition for
  neuroimaging. *PRNI*, 2014.